# Oracle Database 10g: The Power of Globalization Technology

*An Oracle White Paper*
*January 2005*

**ORACLE**®

# Oracle Database 10g: The Power of Globalization Technology

# Oracle Database 10g: The Power of Globalization Technology

## GLOBALIZATION SUPPORT - EXECUTIVE OVERVIEW

Companies today have an unprecedented opportunity to optimize and expand their businesses exponentially. With the advent of the worldwide web Internet and Intranet applications have the possibility of global exposure. With 92% of the world's population being non-native English speaking and the international Internet community being the fastest growing consumers, opportunity knocks. Statistically a customer is 4 times more likely to make a purchase if the content is displayed in their native language. So the challenge and barrier to succeeding in the global market is to have a multilingual enabled application.

How can Oracle help? First you must have a way to be able to easily store, retrieve and update data in any and all languages. Oracle provides full support for Unicode 4.0 the standard for multilingual support. With UTF-8 and UTF-16 support virtually all-contemporary languages and scripts of the world can be easily encoded. This allows customers to develop, deploy, and host multiple languages in a single central database or as part of a grid. Also Oracle offers the flexibility to have all data stored in a Unicode database or incrementally store select columns as a Unicode data type. Another key feature is the ability to present information in the users native localization customs. Things like date, time, currency symbols and delimiters and collation order are handled seamlessly in Oracle. While Oracle's localization support is nearly comprehensive there is the Oracle Locale Builder Utility a graphical tool to do customization for special support.

Typically most web based applications are multi-tier. Oracle provides several database access products for inserting and retrieving Unicode data. Oracle offers support for the most commonly used programming environments such as Java and C/C++. Data is transparently converted between the database and client programs, to ensure that client programs are independent of the database character set.

Putting this all together may require migration of legacy data and applications to a Unicode environment. Oracle provides the Character Set Scanner Utility to be able to proactively detect possible migration problems such as loss or truncation of data and overall impact analysis for the database. Oracle also offers the Globalization Development Kit (GDK) a toolkit that simplifies the development process and

reduces the cost of developing Internet (Java and PL/SQL) applications that will be used to support a global environment.

## THE POWER OF GLOBALIZATION TECHNOLOGY/ INTRODUCTION

So what is the Power of Globalization?  Simply speaking it allows companies to reach and expand their business to "the rest of the world".  Creating application and software products that support the languages and locale customs of the international consumer community.  This is not a trivial task; it requires much more then translating HTML content and imbedded messages.

## GLOBALIZATION SUPPORT IS NATIONAL LANGUAGE SUPPORT AND MORE

Most Oracle customers are familiar with the term National Language Support (NLS) and many have a great deal of experience using it. Every customer that orders a particular release of Oracle gets the same version, same binary file and yet Oracle databases are sold through out the world. How is this possible? National Language Support allows users to store, process, and retrieve data in their native languages and locales. NLS ensures that database utilities, error messages, sort order, date, time, monetary, numeric and calendar conventions automatically adapt to the users native language and locale. This term is somewhat obsolete as Oracle now allows you to do more then simply handle one language for each database instance.

## WHAT IS GLOBALIZATION?

Globalization is the process of developing multilingual applications and software products that can be accessed and run anywhere in the world simultaneously, without modification, rendering content in the native users language and locale preferences. National Language Support is a facilitator for customers to create globalized applications and software. Implementing globalization requires multiple steps. The Process of Globalization can be described as the process of applying Internationalization and Localization.

First let's look at the Internationalization process. In order to have applications that can run anywhere, they must be usable on any language operating system with non-US Keyboards or other country specific hardware. Applications should not have hard-coded dependencies on language strings, and must inter-operate with non-US versions of other products. To be easily translatable applications should isolate language text into separate translatable files. Applications should be code set independent, be able to handle multi-byte characters and handle differences in a distributed environment. To present the correct information the application must be able to detect the users desired locale.

Localization includes translation of the separated file text. Information must then be presented to the user in a manner consistent with the user's local cultural conventions including data formatting, collation, currency, date, time and directionality of text.

**Built in Universal support**

Oracle's National Language Support architecture is implemented with the use of the Oracle NLS Runtime Library (NLSRTL). The NLS Runtime library is not directly exposed to customers but provides a comprehensive suite of API's exposed in SQL, PL/SQL, etc. which allow for proper text and character processing, and culturally and linguistically appropriate data handling. Note, the API's are the same for any locale (universal), but the operations that these API's perform are dependent on the locale settings. This architecture of NLSRTL offers great flexibility in meeting the requirements for software internationalization.

Oracle database applications using SQL, PL/SQL, etc. are built on top of the NLSRTL and the Oracle NLS Architecture, thus inherit the Internationalization model. An application written for example in PL/SQL can be built with one international version that can support any locale. Locale data can be updated without modification to any NLS Runtime Library functions or their usage. Therefore, support for new languages can be added, or the operation of locale-dependent features can be customized by updating the data components. End users and database administrators can then use parameter settings to modify the behavior of the RDBMS to suit their local needs. This provides a flexible architecture that allows language-dependent data to be changed, and data for a new language to be added, without requiring any changes to application code. Also, a single product can support multiple languages.

## WHY GLOBALIZATION?

Go after the largest growing market segment.  92% of the world's population is non-native English speakers.  While the earliest and most fervent users of the Web were of the other 8%, this is quickly changing.  Forrester Research has found that at the moment only 50% of Internet users speak English, and by 2005 it is estimated only a third of internet users will use English for online communication. Making non-native English speakers by far the fastest growing group of Internet users and as we'll discuss next, consumers.  A recent study showed that major web sites in the US turned away almost half the orders coming from outside the country.  Why?  An inability to save, access and retrieve multilingual data such as contact and address information as well as currency issues discourages international customers.  So will non-native users buy from a site even if they can enter multilingual information?  Not likely if the site does not present content in the native users language and locale preferences.  Actually a consumer is 4 times more likely to buy if the content is presented in their native language.

## WHAT'S NEW IN ORACLE GLOBALIZATION SUPPORT?

### Extended Unicode Enablement

What is Unicode? Unicode is a universal encoded character set that allows you to store information from any language using a single character set. Unicode provides a unique code value for every character, regardless of the platform, program, or language.  The Unicode standard has been adopted by many software and hardware vendors, many operating systems and browsers now support Unicode. Unicode is required by modern standards such as XML, Java, JavaScript, LDAP, CORBA 3.0, WML, and it is also compliant to ISO/IEC 10646 standard.

Oracle started supporting Unicode as a database character set in Oracle7. In Oracle9*i* onward, Unicode support has been greatly expanded so that customers can find the right solution for their globalization needs. Oracle Database 10g supports Unicode 4.0, the third and most recent version of the Unicode standard.

#### Unicode Encoding

There are two common ways to encode Unicode 4.0 characters:

- UTF-8 Encoding
- UTF-16 Encoding

#### UTF-8 Encoding

This is the 8-bit encoding of Unicode. It is a variable-width multibyte encoding in which the character codes 0x00 through 0x7F have the same meaning as ASCII. One Unicode character can be 1-byte, 2-bytes, or 3-bytes in this encoding. Generally characters from the European scripts are represented in either 1 or 2

bytes, while characters from most Asian scripts are represented in 3 bytes. Supplementary characters are represented in 4 bytes.

**UTF-16 Encoding**

This is the 16-bit encoding of Unicode in which the character codes 0x0000 through 0x007F have the same meaning as ASCII. One Unicode character is 2-bytes or 4-bytes in this encoding. Characters from both European and most Asian scripts are represented in 2 bytes. Supplementary characters are represented in 4 bytes.

**Unicode Databases**

The Oracle database has the concept of a database character set which specifies the encoding to be used in the SQL CHAR datatypes as well as the metadata such as table names, column names, and SQL statements. A Unicode database must be defined as UTF-8 as the database character set. There are three Oracle character sets that implement the UTF-8 encoding. The first two are designed for ASCII-based platforms while the third one should be used on EBCDIC platforms.

- AL32UTF8

The AL32UTF8 character set is Unicode 4.0 compliant and encodes characters in one to three bytes. Supplementary characters require four bytes. AL32UTF8 is the recommended character set for customers that wish to support UTF-8 in their databases from Oracle Database 10g onward because of the ongoing compliance to the latest Unicode standards.

- UTF8

The UTF8 character set is Unicode 3.0 and CESU-8 compliant and encodes characters in one to three bytes.

- UTFE

The UTFE character set is Unicode 3.0 compliant and should be used as the database character set on EBCDIC platforms to support the UTF-8 encoding.

### New Unicode Datatypes

The SQL NCHAR datatype, which is exclusively Unicode, provides a flexible alternative for multilingual support for customers that need not convert their entire database to Unicode. You can store Unicode characters into columns of these datatypes regardless of the setting of the database character set. The NCHAR datatype was redefined in Oracle9*i* and onward to be a Unicode datatype exclusively. In other words, it stores data in the Unicode encoding only. You can use the SQL NCHAR datatypes in the same way you use the SQL CHAR datatypes.

The encoding used in the SQL NCHAR datatypes is specified as the national character set of the database. You can specify one of the following two Oracle character sets as the national character set:

- AL16UTF16

This is the default character set for SQL NCHAR datatypes. The character set encodes Unicode data in the UTF-16 encoding. Data is counted in 16-bit UTF-16 units.

- UTF8

When UTF8 is specified for SQL NCHAR datatypes, the data stored in the SQL datatypes is in UTF-8, but it is counted as if AL16UTF16 were specified. By default, data is stored in the UTF-16 encoding in the SQL NCHAR datatypes, and the length specified in the NCHAR and NVARCHAR2 columns is always in the number of characters instead of the number of bytes.

### Supplementary Characters

You can extend Unicode to encode more than 1 million characters. These extended characters are called supplementary characters. Supplementary characters are designed to allow representation of characters in future extensions of the Unicode standard. Supplementary characters require 4 bytes in UTF-8 and UTF-16. There is also a reserved area for private use that will never be used by Unicode. This area allows specialized scripts to be stored. As an example perhaps a Star Trek web site wishes to support the language Klingon. Each Klingon character can have a special mapping to the private area using supplementary characters.

### Character Semantics

Character semantics simplifies the task of handling storage requirements and application string manipulation for multibyte strings. Consider character semantics if you are setting up a Unicode database. Unlike say an ASCII character set, UTF-8 and UTF-16 characters are held in multibyte segments. Processing character strings can be much more straightforward especially with variable-width character set such as UTF-8 using character semantics. For example CHAR(10) implies 10 code points or characters of storage. Supported by SQL string functions such as SUBSTR, LENGTH and INSTR.

## Expanded Locale Coverage

Now over 57 languages and 88 countries and territories worldwide and 200 character sets. Through the use of Unicode databases and datatypes, Oracle supports most contemporary languages and scripts. Oracle supports different cultural conventions that are specific to a given geographical location. The default local time format, date format, numeric and monetary conventions are handled based on the local territory setting. By setting different NLS parameters, the database session can utilize different cultural settings. As an example dual currency is supported. In Germany it may be important to display the Deutschmark and the Euro.

Date and Time Formats are also an important consideration to building global applications. The world's various conventions for hour, day, month, and year can be handled in local formats. For example, in the UK, the date is displayed using the format mask 'DD-MON'YYYY' while Japan commonly uses 'YYYY-MON-DD'.

Monetary and Numeric Formats such as currency, credit, and debit symbols also need to be represented in local formats. Radix symbols and thousands separators can be defined by locales. For example, in the US, the decimal separator is a dot ".", while it is a comma "," in France. Therefore, $1,234 could have different meanings in different countries.

Many different calendar systems are in use around the world. Oracle supports seven different calendar systems: Gregorian, Japanese Imperial, ROC Official (Republic of China), Thai Buddha, Persian, English Hijrah, and Arabic Hijrah.

## Overview of Oracle's Sorting Capabilities

Oracle provides linguistic sort capabilities that handle the complex sorting requirements of different languages and cultures. Different languages have different sort orders. What's more, different cultures or countries using the same alphabets may sort words differently. For example, in Danish, the letter Æ is after Z, while Y and Ü are considered to be variants of the same letter. Sort order can be case sensitive or insensitive, and can ignore accents or not. It can also be either phonetic or based on the appearance of the character, such as ordering by the number of strokes or by radicals for East Asian ideographs. Another common sorting issue is when letters are combined. For example, in traditional Spanish, "ch" is a distinct character, which means that the correct order would be: cerveza, Colorado, cheremoya, and so on. This means that the letter "c" cannot be sorted until checking to see if the next letter is an "h". Oracle provides several different types of sort, and can achieve a linguistically correct sort as well as the new multilingual ISO standard (10646) designed to handle many languages at the same time.

**Using Binary Sorts**

Conventionally, when character data is stored, the sort sequence is based on the numeric values of the characters defined by the character encoding scheme. This is called a binary sort. Binary sorts are the fastest type of sort, and produce reasonable results for the English alphabet because the ASCII and EBCDIC standards define the letters A to Z in ascending numeric value. Note, however, that in the ASCII standard, all uppercase letters appear before any lowercase letters. In the EBCDIC standard, the opposite is true: all lowercase letters appear before any uppercase letters. When characters used in other languages are present, a binary sort generally does not produce reasonable results. For example, an ascending ORDER BY query would return the character strings ABC, ABZ, BCD, ÄBC, in the sequence, when the Ä has a higher numeric value than B in the character encoding scheme. For languages using Chinese characters, a binary sort is not linguistically meaningful.

**Using Linguistic Sorts**

To produce a sort sequence that matches the alphabetic sequence of characters, another sort technique must be used that sorts characters independently of their numeric values in the character encoding scheme. This technique is called a linguistic sort. A linguistic sort operates by replacing characters with numeric values that reflect each character's proper linguistic order. These numeric values are found in a table containing major and minor values. Oracle makes two passes when comparing strings. The first pass is to compare the major value of entire string from the major table and the second pass is to compare the minor value from the minor table. Each major table entry contains the Unicode codepoint and major value. Usually, letters with the same appearance will have the same major value. Oracle defines letters with diacritic and case differences for the same major value but different minor values. Oracle offers two kinds of linguistic sort: monolingual, commonly used for European languages; and multilingual, commonly used for Asian languages.

**Using Monolingual Linguistic Sorts**

Oracle offers monolingual linguistic sorts that contain culture-specific sorting order for almost all European languages. Using Multilingual Linguistic Sorts Oracle extends monolingual linguistic sorts so that you can now sort additional languages as part of one sort. This is useful for certain regions or languages that have complex sorting rules or global multilingual databases. Additionally, Oracle still supports all the sort orders defined by the previous releases.

For example a French sort is supported, but the multilingual linguistic sort for French can also be applied by changing the sort order from French to French_M. By doing so, the sorting order will be based on the GENERIC_M sorting order and with the capability to sort secondary level from right to left. Oracle recommends using a multilingual linguistic sort if the tables contain multilingual data. If the tables contain only pure French, for memory usage concern, a French

sort may get better performance. There is a trade-off between extensibility and performance.

For Asian language data or multilingual data, Oracle provides a sorting mechanism based on an ISO standard (ISO14651) and the Unicode 4.0 standard. Multilingual linguistic sorting for Asian languages are implemented in a three-pass fashion based on the number of strokes, Pinyin, or radicals. In addition, handling of canonical equivalence and surrogate codepoint pairs is also implemented with a capacity to define up to 1.1 million codepoints in one sort.

### Using Linguistic Indexes

Using linguistic indices you can provide the sophisticated sorting capabilities of a multilingual sort while achieving sorting performance nearly as good as a binary sort (which offers the best performance).  Function-based index that uses languages other than English can be created.  The index itself does not change the linguistic sort order determined by NLS_SORT. The index simply improves the performance.

### Multiple Linguistic Indexes

If users wish to store character data of multiple languages into one database, they should create multiple linguistic indexes for one column. This approach improves the performance of the linguistic sort for a specific column for multiple languages and is a powerful feature for multilingual databases.

### Case and Accent Insensitive Searching

Operations inside an Oracle database are always sensitive to the case and the accents (diacritics) of the characters. Sometimes you may need to perform case-insensitive or accent-insensitive comparisons and sorts. You could call functions LOWER/UPPER to make SQL statements case-insensitive but this can cause performance degradation.  The new case and accent insensitive searching introduced in Oracle Database 10g is the best of all worlds as it can be applied to any linguistic sort, it won't degrade performance and it allows customers to have the same SQL behaviors without changing existing code.  Just like for linguistic sorts, a function-based index can be built to improve the performance of case-insensitive searches.

### Regular Expression Searching and Replacing

Oracle Regular Expressions provide a simple yet powerful mechanism for rapidly describing patterns and greatly simplifies the way in which you search, extract, format, and otherwise manipulate text in the database.  Traditional regular expression engines were designed to address only English text. However, regular expression implementations can encompass a wide variety of languages with characteristics that are very different from western European text. Oracle's implementation of regular expressions is based on the Unicode Regular Expression

Guidelines. The REGEXP SQL functions work with all character sets that are supported as database character sets and national character sets. Moreover, Oracle enhances the matching capabilities of the POSIX regular expression constructs to handle the unique linguistic requirements of matching multilingual data. Oracle Regular Expression pattern matching is sensitive to the underlying locale defined by the session environment. This affects all aspects of matching including whether it will be case or accent insensitive, whether a character is considered to fall within a range, what collation elements are considered valid, and so on. The engine is also strictly character based, as an example the dot (.) will match a single character in the current character set and never a single byte of the data.

**Date and Time Zones**

Applications that support multi-geographical locales will find comprehensive and precision oriented support for time zones, removing the complexity of doing manual calculations. The datetime data types can store time data with sub-second precision. The datetime data types TSLTZ and TSTZ are time-zone-aware. Datetime values can be specified as local time in a particular region, rather than a particular offset. Using the time zone rules tables for a given region, the time zone offset for a local time is calculated, taking into consideration Daylight Savings time adjustments, and used in further operations.

## Character Set Scanner Utilities

### Migration issues

The Database Character Set Scanner utility provides an assessment of the feasibility and potential issues in migrating an Oracle database to a new database character set. Many customers have discovered and avoided potential migration issues such data truncation, invalid characters and fields that need to be expanded prior to migrating their databases to Unicode using the Character Set Scanner. When migrating to a new database character set, the Export and Import utilities can handle character set conversions from the original database character set to the new database character set. However, character set conversions can sometimes cause data loss or data corruption. For example, if you are migrating from character set A to character set B, the destination character set B should be a superset of character set A. Characters that are not available in character set B will be converted to replacement characters. Another scenario that can cause the loss of data is migrating a database containing data of a different character set from that of the database character set. How can this scenario occur? Users can insert data into the database from another character set if the client `NLS_LANG` character set setting is the same as the database character set. When these settings are the same, Oracle assumes that the data being sent or received is of the same character set, validations and conversions may not be performed. This can lead to two possible data inconsistency problems. One is when a database contains data from another character set but the same codepoints exist in both character sets. The second possibility is having data from mixed character sets inside the database. For example, if the data character set is WE8MSWIN1252, and two separate Windows clients using German and Chinese are both using the `NLS_LANG` character set setting as WE8MSWIN1252, then the database will contain a mixture of German and Simplified Chinese. Obviously for this character set choice i.e. WE8MSWIN1252 the Chinese characters are not expected or supported.

### Anticipating Migration Issues

The Scanner checks all character data in the database and tests for the effects and problems of changing the character set encoding. At the end of the scan, it generates a summary report of the database scan. This report provides estimates of the amount of work required to convert the database to a new character set. The Scanner reads the character data and tests for the following conditions on each data cell:

- Do character codes of the data cells change when converted to the new character set?

- Can the data cells be successfully converted to the new character set?

- Will the post-conversion data fit into the current column size?

The Scanner reads and tests for data in CHAR, VARCHAR2, LONG, CLOB, NCHAR, NVARCHAR2, and NCLOB columns only. The Scanner does not perform post-conversion column size testing for LONG, CLOB, and NCLOB columns.

### Scanning HTML and Plain Text Files

Loading files to the database safely requires knowledge of the potential conversion that may take place. The Language and Character Set File Scanner (LCSSCAN) Utility introduced in Oracle Database 10g removes the guesswork of identifying the encoding of your source file. LCSSCAN is a statistically based client utility for determining the character set and language for unspecified file text.

## The Oracle Locale Builder

The Oracle Locale Builder allows users to customize virtually any type of locale definition simply and safely through an intuitive graphical interface. The Oracle Locale Builder offers an easy and efficient way to access and define NLS locale data definitions. It provides a graphical user interface through which users can easily view, modify, and define various locale-specific data. It extracts data from the definition files (`*.nlt` and `*.nlb`) and presents them in a readable format, to allow processing of the information without worrying about the specific definition formats used in these files. Users can choose to have the resulting output file be in either text format (.nlt) or the default binary format (.nlb).

Oracle Locale Builder can be accessed either locally by running it as a local application or remotely with a standard web browser.

The Oracle Locale Builder handles 4 types of locale definitions:

- **Territories** including calendar convention, date and time formats, number and monetary systems

- **Languages** including local month and day names, and writing direction

- **Character Sets** including character set type, character mappings, character set charts, and classifications

- Linguistic Sorts including linguistic sort order, and special collation rules

Locale Builder offers 2 approaches to adding new Locale definitions or customizing existing definitions. To do this essentially you build upon an inherited setting and modify it to add additional properties. For example perhaps you wish to build a hybrid language French American. With the Locale Builder you could choose French. Modify month and day names or character rules to create a Unique French American language. As opposed to creating a whole new language you may want to customize a new territory. After defining a new territory for a given language you can then modify date, time, monetary and numeric formats without affecting the Language characteristics. Numeric formats include decimal separator, how numbers are rounded and what measurement standard is used such as the metric

**Transportable NLB Data**

**Starting in Oracle Database 10g, NLB files are transportable. After using Oracle Locale Builder to customize a locale the NLB files that are generated can be transported to another platform by, for example, FTP. The transported NLB files can be used the same way as the NLB files that were generated on the original platform.**

system.  Monetary formats can be changed such as currency symbol, decimal and group separators grouping and precision, and credit and debit symbols.

You can create new character sets and extend an existing encoded character set definition by using User-defined characters (UDC) to encode special characters representing:

- Proper names

- Historical Han characters that are not defined in an existing character set standard

- Vendor-specific characters

- New symbols or characters you define

User-defined characters are typically supported within such as Unicode and East Asian character sets.  Character sets that support User-defined characters have at least one range of reserved codepoints for use as user-defined characters. For example, Japanese Shift JIS preserves 1880 codepoints for user-defined characters.

Linguistic Sort order can be changed in many ways.  Code point ranges can be reconfigured so for example numbers sort after letters or you can change the order of individual code points.  You can change the sort for all characters containing a particular diacritic or change one diacritic character at a time.  Canonical equivalence can be configured or turned off as needed.

### Globalization Development Kit

The Globalization Development Kit (GDK) is a set of Java APIs that provide Oracle application developers with the framework to develop globalized Internet applications using the best globalization practices and features designed by Oracle. The GDK complements the existing globalization features in Java; it provides the synchronization of locale behaviors between the middle tier Java application and the Oracle database server.

The functionalities offered by GDK for Java can be divided into two categories:

- The GDK application framework for J2EE provides the globalization framework for building J2EE-based Internet application. The framework encapsulates the complexity of globalization programming, such as determining user locale, maintaining locale persistency, and processing locale information. It consists of a set of Java classes through which applications can gain access to the framework. These associated Java classes enable applications to code against the framework so that globalization behaviors can be extended declaratively.

- The GDK Java API offers development support in Java applications and provides consistent globalization operations as provided in Oracle database servers. The API is accessible and is independent of the GDK

framework. Therefore standalone Java applications and J2EE applications that are not based on the GDK framework are able to access the individual features offered by the Java API. The features provided in the Java API include data and number formatting, sorting, and handling character sets in the same way as the Oracle Database. Language and Character Set detection technology is also available via the GDK Java API.

The GDK is certified with both Oracle9*i* and Oracle Database 10*g* databases running JDK version 1.3 or above.

## Building a Central Multilingual Server

### Architectural Diagram

An e-business runs on one globalized computer system. Everybody is connected. And all the information is shared in one place.

### Taking a Central Server Approach

Here is what Oracle faced prior to creating a single point solution and saving over a billion dollars in the process. Does this sound like your organization?  Information was scattered across hundreds of separate databases and that was the problem! Each organizations, marketing, sales, service, etc. had its own computer system. Each computer system had its own database. Oracle had hundreds of databases around the world. Data was so fragmented, it was difficult for people to find the information they needed to do their jobs. Separate databases also made it difficult to share information between the organizations. And if groups can't share information, they don't cooperate. So, marketing didn't cooperate with sales. Germany didn't cooperate with France. And the lack of cooperation led to duplication of effort and inefficiency. To eliminate this inefficiency, Oracle had to make information easier to find and easier to share. But how? The solution was quite simple become an e-business.  But what is an e-business anyway? It's all about the Internet and globalization. An e-business uses a global network and a global database to integrate all aspects of doing business. Every business function: marketing, sales, supply chain, manufacturing, customer service, accounting, human resources, everything uses the same global network and the same global database. An e-business runs on one unified computer system. Everybody is connected. And all the information is in one place.  While conceptually simple, this single, unified database approach required fundamental changes to application software. Most companies whether doing B2B, B2C or hosting face the challenge of handling

multilingual data and presenting information in their users native languages and locales. Middle tier software must be able to support web traffic coming from anywhere in the world, and identify the customers language and locale preferences, and create the proper translated pages.  If we are to consolidate servers and data to a unary approach then all components must be multilingual ready.

**Implementing a Unicode Solution in the Database**

The first step in building a central multilingual server is to have a way to be able to easily store, retrieve and update data in any and all languages. With UTF-8 and UTF-16 support virtually all contemporary languages and scripts of the world can be easily encoded.  This allows customers to develop, deploy, and host multiple languages in a single central database.  Multilingual information can be shared and yet each user can be served in their own locale preferences.

**Implementing a Single Multilingual Application Server with the GDK**

Often applications start out to support a single language and as the system grows the requirements for multilingual support evolve.  Creating multiple versions of the same application that support an individual language is an option but it is resource and maintenance intensive. Supporting multiple locales simultaneously from a single binary and a single code base has many benefits including reducing the cost of hardware, maintenance and faster support of additional locales. Unicode lets your applications support multiple languages simultaneously. And the single multilingual server approach provides full Unicode support in all tiers of the application architecture.

Customers shy away from this solution because of their lack of knowledge on internationalizing applications.  Let's face it developers have enough challenges creating object oriented applications and applying the company business logic.  Oracle's Globalization Development Kit (GDK) simplifies the development process and reduces the cost of developing Internet applications that will be used to support a global environment.

**Manage your applications with the Oracle Application Server**

Just as the Oracle database manages all your data, Oracle Application Server (runs all your applications. For Internet and intranet applications, Oracle Application Server offers the most innovative and comprehensive set of middle-tier services.  While offering unique capabilities to quickly develop web applications a second key feature is the tight integration of Oracle Application Server with Oracle database technologies.  Each Oracle Application Server development environment provides programming methods that minimize the complexities of inserting and retrieving Unicode data regardless of the client character set.  For example let's look at JSP, Java Servlets, and PSPs.

Oracle Application Server implements the standard Java Servlet API that allows the deployment of Java Servlets. It also implements the JSP compiler according to the

standard Java ServerPages specification that allows users to compile standard JSPs to Java Servlets. As a result, applications can fully utilize the internationalization support provided in the Java (JDK), Java Servlet and JSP technologies. Oracle JDBC drivers transparently convert the data from the database character set if needed to UTF-16, for all Char types, as well as LONG and CLOB. As a result of this transparent conversion, JSPs and Java Servlets calling Oracle JDBC drivers may bind and define database columns with Java strings, and fetch data into Java strings from the result set of a SQL execution.

PSP and PL/SQL Web Toolkit Oracle iAS provides a Web gateway that allows PL/SQL stored procedures to generate dynamic web content and deliver it to the client browser in the same way as Java Servlets. Oracle provides an API called the PL/SQL Web Toolkit for the development of Internet applications in PL/SQL stored procedures. The API provides methods for you to format web pages in PL/SQL and send them out to the client browser via the Web gateway. In addition to the Web Toolkit, you may also use the SQL functions, such as SUBSTRING(), TO_DATE() and LENGTH(), provided by Oracle to manipulate strings. All string variables, such as VARCHAR and CHAR strings, and the SQL functions used in the PL/SQL stored procedures operate on Unicode.

PSPs are HTML pages with embedded PL/SQL code. PSP relates to PL/SQL stored procedure in the same way as JSP relates to Java Servlet. Oracle provides a PSP compiler to compile PSPs into a PL/SQL stored procedures and load them into the database. When using PSP and PL/SQL stored procedures in an application, direct access to UTF-8 and UTF-16 data comes automatically by using SQL or PL/SQL from inside the database.

## Summary

### Flexible

With Oracle globalization support customers can create a locale environment that meets their business needs, from basic monolingual, native language support to centralized multilingual databases. Localization parameter settings can be specified from the client session, server, or explicitly within SQL functions. Customization features are also available to help tailor the environment. Customers needing to move to a Unicode solution can choose between using a Unicode database or the Unicode datatypes to best suit their business needs.

### Compatible

Support for 200 different national, multinational, and vendor-specific character sets, provides consistency and interoperability with existing data. The most popular single-byte, multi-byte, and fixed-width encodings used in the industry today; including UTF-8 and UTF-16 for full Unicode 4.0 support. For customers needing to move to a Unicode solution will find, the Character Set Scanner utility can help identify compatibility issues to assure a smooth migration. Oracle further supports deployment of multi-tier multilingual applications by automatically and transparently performing any necessary character set conversions on the database. The Globalization Development Kit (GDK) contains a set of Java APIs that provide Oracle application developers with the framework to develop globalized Internet applications using the best globalization practices and features designed by Oracle GDK complements the existing globalization features in J2EE. Although the J2EE platform already provides a strong foundation for building globalized applications, its globalization functionalities and behaviors can be quite different from Oracle's functionalities. GDK provides synchronization of locale-sensitive behaviors between the middle-tier Java application and the database server. In addition the GDK brings some of the best database globalization features to the middle tier. The GDK also contains a suite of PL/SQL packages that provide additional globalization functionalities for applications written in PL/SQL.

### Integrated

Oracle makes use of a fully internationalized single-binary model, which ensures that the same Oracle release can service users all around the globe. UI translations are also included in many different languages. Plan to support multilingual applications and a Unicode database? The Oracle Application Server offers the most innovative and comprehensive set of middle-tier services. All access programming interfaces to Oracle are enabled for both UTF-16 and UTF-8, thus providing excellent native integration for applications written in these Unicode forms. Oracle also offers the Globalization Development Kit (GDK) a toolkit that removes the complexity of internationalizing an application and let's developers focus on the company business logic.

# ORACLE®

**The Power of Globalization**
**January 2005**
**Author: Barry Trute**
**Contributing Authors:**

**Oracle Corporation**
**World Headquarters**
**500 Oracle Parkway**
**Redwood Shores, CA 94065**
**U.S.A.**

**Worldwide Inquiries:**
**Phone: +1.650.506.7000**
**Fax: +1.650.506.7200**
**www.oracle.com**